

Adaptive Approximation of Functions with Discontinuities

Licia Lenarduzzi and Robert Schaback

Version of Nov. 09, 2015

Abstract: One of the basic principles of Approximation Theory is that the quality of approximations increase with the smoothness of the function to be approximated. Functions that are smooth in certain subdomains will have good approximations in those subdomains, and these *sub-approximations* can possibly be calculated efficiently in parallel, as long as the subdomains do not overlap. This paper proposes a class of algorithms that first calculate sub-approximations on non-overlapping subdomains, then extend the subdomains as much as possible and finally produce a global solution on the given domain by letting the subdomains fill the whole domain. Consequently, there will be no Gibbs phenomenon along the boundaries of the subdomains. Throughout, the algorithm works for fixed scattered input data of the function itself, not on spectral data, and it does not resample.

Key words: Kernels, classification, localized approximation, adaptivity, scattered data

AMS classification: 65D05, 62H30, 68T05

1 Introduction

Assume that a large set $\{(\mathbf{x}_i, f_i), i = 1, \dots, N\}$ of data is given, where the points \mathbf{x}_i are scattered in \mathbb{R}^d and form a set X . We want to find a function u that recovers the data on a domain Ω containing the points, i.e.

$$\begin{aligned} u &: \Omega \rightarrow \mathbb{R}, \\ u(\mathbf{x}_i) &\approx f_i, i = 1, \dots, N. \end{aligned}$$

We are particularly interested in situations where the data have smooth interpolants in certain non-overlapping subdomains Ω_j , but not globally. The reason may be that there are discontinuities in the function itself or its derivatives. Thus a major goal is to identify subdomains $\Omega_j \subseteq \Omega$, $1 \leq j \leq J$ and smooth functions u_j , $1 \leq j \leq J$ such that

$$\begin{aligned} u_j &: \Omega_j \rightarrow \mathbb{R}, \\ u_j(\mathbf{x}_i) &\approx f_i \text{ for all } \mathbf{x}_i \in X \cap \Omega_j. \end{aligned}$$

The solution to the problem is piecewise defined as

$$u(\mathbf{x}) := u_j(\mathbf{x}) \text{ for all } \mathbf{x} \in \Omega_j, 1 \leq j \leq J.$$

Our motivation is the well-known fact that errors and convergence rates in Approximation Theory always improve with increasing smoothness. Thus on each subdomain we expect to get rather small errors, much smaller than if the problem would have been treated globally, where the non-smoothness is a serious limiting effect.

From the viewpoint of Machine Learning [3, 8, 9] this is a mixture of classification and regression. The domain points have to be classified in such a way that on each class there is a good regression model. The given training data are used for both classification and regression, but in this case the classification is dependent on the regression, and the regression is dependent on the classification.

Furthermore, there is a serious amount of geometry hidden behind the problem. The subdomains should be connected, their interiors should be disjoint, and the union of their closures should fill the domain completely. This is why a black-box machine learning approach is not pursued here. Instead, Geometry and Approximation Theory play a dominant part. For the same reason, we avoid to calculate edges or fault lines first, followed by local approximations later. The approximation properties should determine the domains and their boundaries, not the other way round.

In particular, *localized approximation* will combine Geometry and Approximation Theory and provide a central tool, together with *adaptivity*. The basic idea is that in the interior of each subdomain, far away from its boundary, there should be a good and simple approximation to the data at each data point from the data of its neighbors.

2 An Adaptive Algorithm

Localized approximation will be used as the first phase of an *adaptive algorithm*, constructing disjoint localized subsets of the data that allow good and simple *local approximations*. Thus this “localization” phase produces a subset $X^g \subseteq X$ of “good” data points that is the union of disjoint sets X_1^g, \dots, X_J^g consisting of data points that allow good approximations $u_j^g \in U$, $1 \leq j \leq J$ using only the data points in X_j^g . In some sense, this is a rough classification already, but only of data points.

The goal of the second phase is to reduce the number of unclassified points by enlarging the sets of classified points. It is tacitly assumed that the final number

of subdomains is already obtained by the number J of classes of “good” points after the first phase. The “blow-up” of the sets X_j^g should maintain *locality* by adding neighboring data points first, and adding them only if the local approximation u_j^g does not lose too much quality after adding that point and changing the approximation.

The second phase usually leaves a small number of “unsure” points that could not be clearly classified by blowing up the classified sets. While the blow-up phase focuses on each single set X_j^g in turn and tries to extend it by looking at all “unsure” points for good extension candidates, the third phase works the other way round. It focuses on each single “unsure” point \mathbf{x}_i in turn and looks at all sets X_j^g and the local approximations u_j on these, and assigns the point \mathbf{x}_j to one of the sets X_j^g so that $u_j(\mathbf{x}_i)$ is closest to $f(\mathbf{x}_i)$. It is a “final assignment” phase that should classify all data points and it should produce the final sets $X_j^f \supseteq X_j^g$ of data points. The sets X_j^f should be disjoint and their union should be X .

After phase 3, each local approximation $u_j^f \in U$ is based on the points in X_j^f only, but there still are no well-defined subdomains $\Omega_j \supseteq X_j^f$ as domains of u_j^f . Thus the determination of subdomain boundaries from a classification of data points could be the task of a fourth phase. It could, for instance, be handled by any machine learning program that uses the classification as training data and classifies each given point \mathbf{x} accordingly. But this paper does not implement a fourth phase, being satisfied if each approximation u_j^f is good on each set X_j^f , and much better than any global approximation $u^* \in U$ to all data.

3 Implementation

The above description of a three-phase algorithm allows a large variation of different implementations that compete for efficiency and accuracy. We shall describe a basic implementation together with certain minor variants, and provide numerical examples demonstrating that the overall strategy works fine.

We work on the unit square of \mathbb{R}^2 for simplicity and take a trial space U spanned by translates of a fixed positive definite radial kernel K . In our examples, K may be a Gaussian or an inverse multiquadric. For details on kernels, readers are referred to standard texts [2, 10, 7, 4], for example. When working on finite subsets of data points, we shall only use the translates with respect to this subset. Since the kernel K is fixed, also the Hilbert space H is fixed in which the kernel is reproducing, and we can evaluate the norm $\|\cdot\|_K$ of trial functions

cheaply and exactly.

To implement locality, we assume that we have a computationally cheap method that allows to calculate for each $\mathbf{x} \in \mathbb{R}^2$ its n nearest neighbors from X . This can, for instance, be done via a range query after an initialization of a kd-tree data structure [1].

3.1 Phase 1: Localization

This is carried out by a first step picking all data points with good localized approximation properties, followed by a second step splitting the set X^g of good points into J disjoint sets X_j^g .

3.1.1 Good Data Points

We assume that the global fill distance

$$h(X, \Omega) := \sup_{\mathbf{y} \in \Omega} \min_{\mathbf{x}_k \in X} \|\mathbf{y} - \mathbf{x}_k\|_2$$

of the full set of data points with respect to the full domain Ω is roughly the same as the local fill distances $h(X_j^f, \Omega_j)$ of the final splitting.

The basic idea is to loop over all N data points of X and to calculate for each data point \mathbf{x}_i , $1 \leq i \leq N$ a number σ_i that is a reliable indicator for the quality of *localized approximation*. Using a threshold σ , this allows to determine the set $X^g \subseteq X$ of “good” data points, without splitting it into subsets.

There are many ways to do this. The implementation of this paper fixes a number n of neighbors and loops over all N data points to calculate for each data point \mathbf{x}_i , $1 \leq i \leq N$

1. the set N_i of their n nearest neighbors from X ,
2. the kernel-based interpolant s_i of the data $(\mathbf{x}_k, f(x_k))$ for all n neighboring data points $\mathbf{x}_k \in N_i$,
3. the norm $\sigma_i := \|s_i\|_K$.

This loop can be executed with roughly $\mathcal{O}(Nn^3)$ complexity and $\mathcal{O}(N + n^3)$ storage, and with easy parallelization, if necessary at all. A similar indicator would be the error obtained when predicting $f(\mathbf{x}_i)$ from the values at the n neighboring points.

Practical experience shows that the numbers σ_i are good indicators of locality, because adding outliers to a good interpolant usually increases the error norm dramatically. Many of the σ_i can be expected to be small, and thus the threshold

$$\sigma_i < 2M_\sigma$$

will be used to determine “good” points within the next splitting step, see Section 3.1.2, where M_σ is the median of all σ_i . This is illustrated for a data set by Figure 1: it represents, in base loglog scale, the sorted $\{\sigma_i\}$ and the constant line relevant to the value of the threshold.

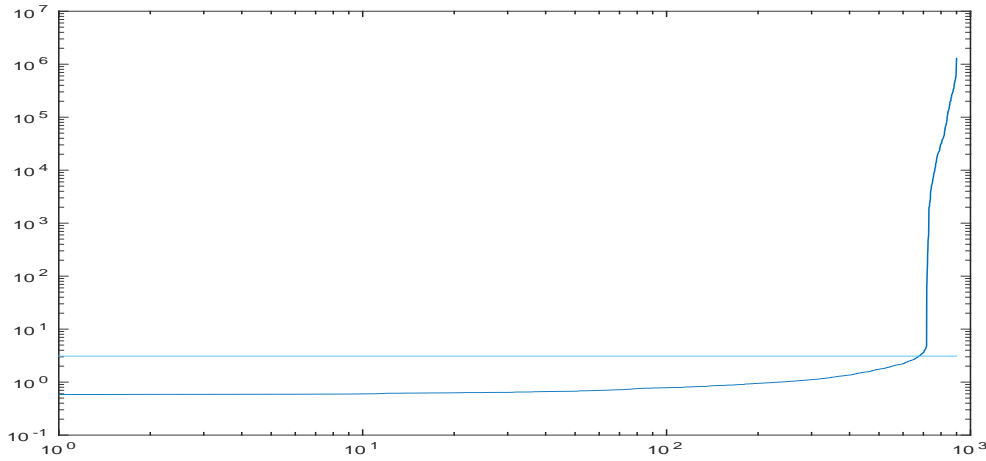


Figure 1: Loglog: sorted $\{\sigma_i\}$ and threshold

3.1.2 Splitting

The set X^g of points with good localization must now be split into J disjoint subsets of points that are close to each other.

We assume that the inner boundaries of the subdomains are everywhere clearly determined by large values of σ_i .

The implementation of this paper accomplishes the splitting by a variation of Kruskal’s algorithm [5] for calculating minimal spanning trees in graphs.

The Kruskal algorithm sorts the edges by increasing weight and starts with an output graph that has no edges and no vertices. When running, it keeps a number of disconnected graphs as the output graph. It gradually adds edges with

increasing weight that either connect two previously disconnected graphs or add an edge to an existing component or define a new connected component by that single edge.

In the current implementation the edges that connect each x_i with its $n - 1$ nearest neighbors are collected in an edge list. The edge list is sorted by increasing length of the edges and then, by $n \mid X \mid$ comparisons, many repetitions of edges are removed, and these are all repetitions if any two different edges have different length.

Then the thresholding of the $\{\sigma_i\}$ by

$$\sigma_i < 2M_\sigma$$

is executed, and it is known which points are good and which are bad.

All edges with one or two bad end points are removed from the edge list with cost $n \mid X \mid$.

After the spanning tree algorithm is run, the list of the points of each tree is intersected with itself to avoid eventual repetitions that are left. At the end, each connected component is associated to its tree in exactly one way.

In rare cases, the splitting step may return only one tree, but these cases are detected and repaired easily.

3.2 Phase 2: Blow-up

This is also an adaptive iterative process. It reduces the set $X'' := X \setminus \bigcup_{j=1}^J X_j^g$ of “unsure” data points gradually, moving points from X'' to one of the nearest sets X_j^g . In order to deal with easy cases first, the points \mathbf{x}_i in X'' are sorted by their locality quality such that points with better localization come first. We also assume that for each point $\mathbf{x}_i \in X''$ we know its distance to all sets X_j^g , and we shall update this distance during the algorithm, when the sets X'' and X_j^g change. We also use the distances to the sets $X_j^{g,0}$ that are the output of the localization phase and serve as a start-up for the sets X_j^g .

In an outer loop we run over all points $\mathbf{x}_i \in X''$ with decreasing quality of local approximation. In our implementation, this means increasing values of σ_i . The inner loop runs over the m sets X_j^g to which \mathbf{x}_i has shortest distance. In most cases, and in particular in \mathbb{R}^2 , it will suffice to take $m = 2$. The basic idea is to find the nearby set X_j^g of “good” points for which the addition of \mathbf{x}_i does least damage to the local approximation quality.

Our implementation of the inner loop over m neighboring sets X_j^g works as follows. In $X_j^{g,0}$, the point \mathbf{y}_j with shortest distance to \mathbf{x}_i is picked, and its n nearest neighbors in $X_j^{g,0}$ are taken, forming a set Y_j^g . On this set, the data interpolant s_j^g is calculated, and then the number $\sigma_j^g := \|s_j^g\|_K$ measures the local approximation quality near the point \mathbf{y}_j if only “good” points are used. Then the “unsure” point \mathbf{x}_i is taken into account by forming a set Y_j^u of points consisting of \mathbf{x}_i and the up to $n - 1$ nearest neighbors to \mathbf{x}_i from X_j^g . On this set, the data interpolant s_j^u is calculated, and the number $\sigma_j^u := \|s_j^u\|_K$ measures the local approximation quality if the “unsure” point \mathbf{x}_i is added to X_j^g . The inner loop ends by maintaining the minimum of quotients σ_j^u/σ_j^g over all nearby sets X_j^g checked by the loop. These quotients are used to indicate how much the local approximation quality would degrade if \mathbf{x}_i would be added to X_j^g . Note that this strategy maintains locality by focusing on “good” nearest neighbors of either \mathbf{x}_i or \mathbf{y}_j . By using the fixed sets $X_j^{g,0}$ instead of the growing sets X_j^g , the algorithm does not rely heavily on the newly added points.

An illustration is attached to Example 1 in the next section; there the point \mathbf{y}_1 and the sets X_1^g and Y_1^u associated to a point \mathbf{x}_i will be shown.

After the inner loop, if the closest set to \mathbf{x}_i among all sets X_k^g is X_j^g and σ_j^u/σ_j^g is less than σ_k^u/σ_k^g for $k \neq j$, then \mathbf{x}_i is moved from X^u to X_j^g . If the closest set to \mathbf{x}_i is X_j^g but if it is not true that σ_j^u/σ_j^g is less than σ_k^u/σ_k^g for $k \neq j$, then \mathbf{x}_i remains “unsure”. The “unsure” points are those that seriously degrade the local approximation quality of all nearby sets of “good” points.

3.3 Phase 3: Final Assignment

The assignment of a point $\mathbf{x}_i \in X^u$ to a set X_j^g is done on the basis of how well the function value $f(\mathbf{x}_i)$ is predicted by $u_j(\mathbf{x}_i)$. We loop over all points $\mathbf{x}_i \in X^u$ and first determine two sets X_j^g and X_k^g to which \mathbf{x}_i has shortest distance. This is done in order to make sure that \mathbf{x}_i is not assigned to a far-away X_j^g . We then could add \mathbf{x}_i to X_j^g if $|f(\mathbf{x}_i) - u_j(\mathbf{x}_i)| \leq |f(\mathbf{x}_i) - u_k(\mathbf{x}_i)|$, otherwise to X_k^g , but in case that we have more than one unsure point, we want to make sure that under all unsure

points, \mathbf{x}_i fits better into X_j^g than into X_k^g . Therefore we calculate

$$\begin{aligned} d_j(\mathbf{x}_i) &:= |f(\mathbf{x}_i) - u_j(\mathbf{x}_i)| \\ \mu_j &:= \min_{\mathbf{x}_i \in X^u} d_j(\mathbf{x}_i) \\ M_j &:= \max_{\mathbf{x}_i \in X^u} d_j(\mathbf{x}_i) \\ D_j(\mathbf{x}_i) &:= \frac{d_j(\mathbf{x}_i) - \mu_j}{M_j - \mu_j} \end{aligned}$$

for all j and i beforehand, and assign \mathbf{x}_i to X_j^g if $D_j(\mathbf{x}_i) \leq D_k(\mathbf{x}_i)$, otherwise to X_k^g .

4 Examples

Some test functions are considered now, each of which is smooth on $J = 2$ subdomains of Ω . The algorithm constructs X_1^g and X_2^g with $X_1^g \cup X_2^g = X$.

Concerning the error of approximation of u , we separate what happens away from the boundaries of Ω_j from what happens globally on $[0, 1]^2$. This is due to the fact that standard domain boundaries, even without any domain splittings, let the approximation quality decrease near the boundaries.

To be more precise, let Ω_{safe} be the union of the circles of radius

$$q := \min_{1 \leq i < j \leq N} \|\mathbf{x}_i - \mathbf{x}_j\|_2,$$

the separation distance of the data sites, centered at those points of X_j^g , $j = 1, 2$ such that the centered circles of radius $2q$ do not contain points of X_k^g with $k \neq j$. We then evaluate

$$L_\infty^{safe}(u) := \|u - f\|_{\infty, \Omega_{safe} \cap [0, 1]^2}$$

and

$$L_\infty(u) := \|u - f\|_{\infty, [0, 1]^2}.$$

The chosen kernel for calculating the local kernel-based interpolants is the inverse multiquadric kernel $\phi(r) = (1 + 2r^2/\delta^2)^{-1/2}$ with parameter $\delta = 0.35$.

In all cases, $N = 900$ data locations are mildly scattered on a domain Ω that extends $[0, 1]^2$ a little, with $q = 0.04$. We shall restrict to $[0, 1]^2$ to evaluate the subapproximant, calculated by the basis in the Newton form. Such a basis is much more stable than the standard basis, see [6]. The error is computed on a grid with step 0.01.

Example 1. The function

$$f_1(x, y) := \log(|x - (0.2 \sin(2\pi y) + 0.5)| + 0.5),$$

has a derivative discontinuity across the curve $x = 0.2 \sin(2\pi y) + 0.5$. We get

$$L_\infty^{sae}(u) = 1.6 \cdot 10^{-5}, \quad L_\infty(u) = 6.0 \cdot 10^{-2}.$$

For comparison, the errors of the global interpolant are

$$L_\infty^{sae}(u^\star) = 1.1 \cdot 10^{-1}, \quad L_\infty(u^\star) = 1.1 \cdot 10^{-1}.$$

The classification turns out to be correct. 890 out of 900 data points are correctly classified as output of phase 3.2, and then phase 3.3 completes the classification.

Figure 2 shows the points of X_1^f as dotted and those of X_2^f as crossed. The points both dotted and circled of X_1^f , respectively the points both crossed and circled of X_2^f , are the result of the splitting (Section 3.1.2), while the points dotted only, respectively crossed only, are those added by the blow-up phase (Section 3.2). The points squared are the result of the final assignment phase (Section 3.3). The true splitting line is traced too. The convention of the marker types will be used in the next examples as well.

The function u is defined as u_1^f where the subdomain Ω_1 is determined and as u_2^f on Ω_2 .

The actual error $L_\infty(u) = 6.0 \cdot 10^{-2}$ is not much affected if we omit Phase 3 and and ignore the remaining 10 “unsure” points after the blow-up phase. A similar effect is observed for the other examples to follow.

A zoomed area of Ω is considered in Figure 3. The details are related to an iteration of the blow-up phase, where the “unsure” point \mathbf{x}_i (both squared and starred) is currently examined. Points of $X_2^{g,0}$ are shown as crosses. At the current iteration, the dots are points inserted in X_1^g up to now, those belonging to $X_1^{g,0}$ bold dotted, while the points inserted in X_2^g up to now are omitted in this illustration. The points squared are those of Y_1^u , while the points as diamonds are those of Y_1^g . The point \mathbf{y}_1 is both written as diamond and star.

Example 2. The function

$$f_2(x, y) := \begin{cases} f_1(x, y) & \text{if } x \leq 0.2 \sin(2\pi y) + 0.5 \\ f_1(x, y) + 0.01 & \text{if } x > 0.2 \sin(2\pi y) + 0.5 \end{cases}$$

has a discontinuity across the curve $x = 0.2 \sin(2\pi y) + 0.5$.

We get

$$L_{\infty}^{safe}(u) = 1.6 \cdot 10^{-5}, L_{\infty}(u) = 6.0 \cdot 10^{-2}.$$

For comparison, the errors of the global interpolant are

$$L_{\infty}^{safe}(u^*) = 1.3 \cdot 10^{-1}, L_{\infty}(u^*) = 1.3 \cdot 10^{-1}.$$

The classification turns out to be correct. 888 out of 900 data points are classified correctly as output of phase 3.2, and phase 3.3 completes the classification for the remaining 12 points. It might be that u_j^f is more accurate on the safe zone, and also globally.

Example 3. The function

$$f_3(x, y) := \arctan(10^3(\sqrt{(x+0.05)^2 + (y+0.05)^2} - 0.7)) \quad (1)$$

is regular but has a steep gradient. Our algorithm yields

$$L_{\infty}^{safe}(u) = 9.0 \cdot 10^{-2} \text{ and } L_{\infty}(u) = 2.67 \cdot 10^0,$$

while for the global interpolant we get

$$L_{\infty}^{safe}(u^*) = 2.31 \cdot 10^0, L_{\infty}(u^*) = 3.26 \cdot 10^0.$$

Figure 4 shows the points of X_1^f as dotted and those of X_2^f as crossed; X_1^f and X_2^f stay at the opposite sides of the mid range line $f(x, y) = 0$.

Example 4. The function

$$f_4(x, y) := ((x - 0.5)^2 + (y - 0.5)^2)^{0.35} + 0.05 * (x - 0.5)_+^0 \quad (2)$$

has a jump on the line $x = 0.5$ and a derivative singularity on it at $(0.5, 0.5)$. It has rather a steep gradient too. One data point close to the singularity is not classified correctly. We get

$$L_{\infty}^{safe}(u) = 9.9 \cdot 10^{-4} \text{ and } L_{\infty}(u) = 7.3 \cdot 10^{-2},$$

while the global interpolant u^* has

$$L_{\infty}^{safe}(u^*) = 1.5 \cdot 10^{-2} \text{ and } L_{\infty}(u^*) = 8.5 \cdot 10^{-2}.$$

Figure 5 shows the points of X_1^f as dotted and those of X_2^g as crossed.

All examples show that the transition from a global to a properly segmented problem decreases the achievable error considerably. But the computational cost is serious, and it might be more efficient to implement a multiscale strategy that works on coarse data first, does the splitting of the domain coarsely, and then refines the solution on more data, without recalculating everything on the finer

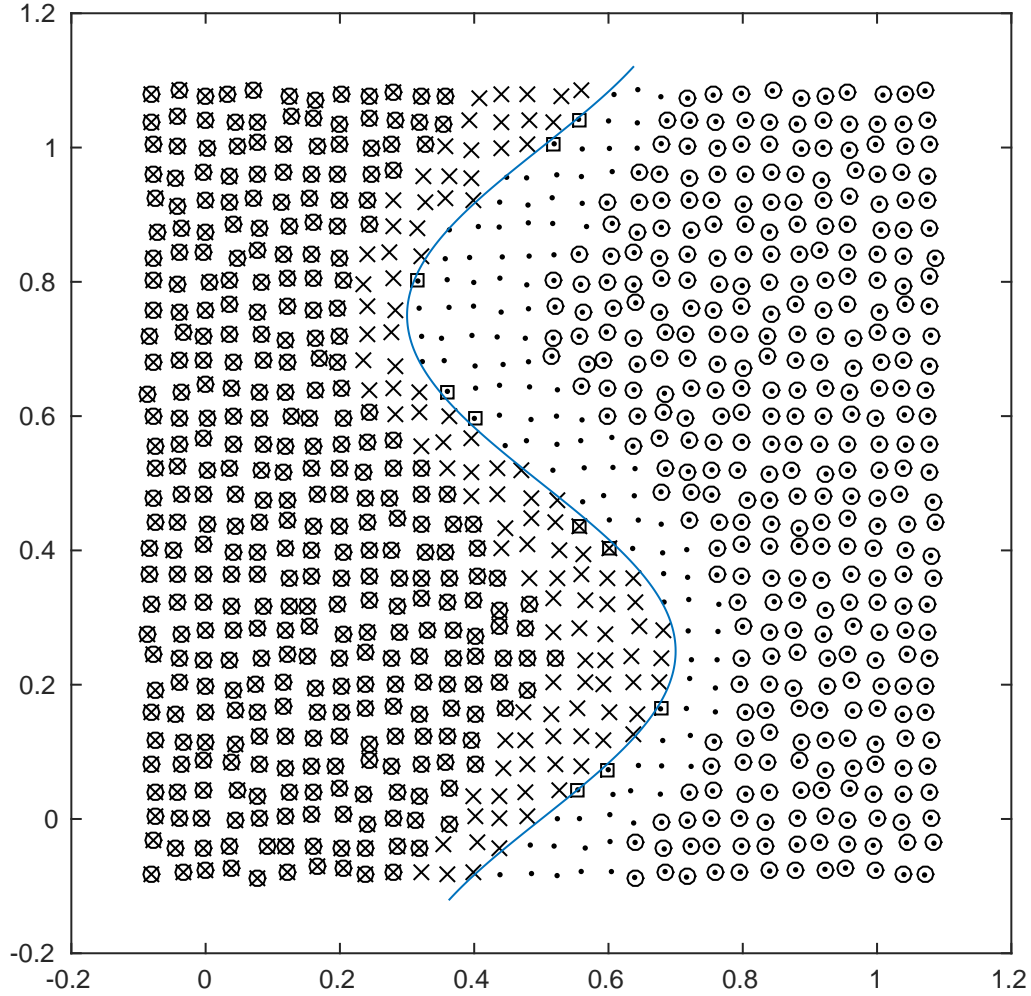


Figure 2: Example 1: class 1 as dots, class 2 as crosses

References

- [1] J.L. Bentley. Multidimensional binary search trees used for associative searching. *Communications of the ACM*, 18:509–517, 1975.
- [2] M.D. Buhmann. *Radial Basis Functions, Theory and Implementations*. Cambridge University Press, 2003.
- [3] N. Cristianini and J. Shawe-Taylor. *An introduction to support vector machines and other kernel-based learning methods*. Cambridge University Press, Cambridge, 2000.
- [4] G. Fasshauer and M. McCourt. *Kernel-based Approximation Methods using MATLAB*, volume 19 of *Interdisciplinary Mathematical Sciences*. World Scientific, Singapore, 2015.
- [5] J.B. Kruskal. On the shortest spanning subtree of a graph and the traveling salesman problem. *Proceedings of the American Mathematical Society*, 7:48–50, 1956.
- [6] St. Müller and R. Schaback. A Newton basis for kernel spaces. *Journal of Approximation Theory*, 161:645–655, 2009. doi:10.1016/j.jat.2008.10.014.
- [7] R. Schaback and H. Wendland. Kernel techniques: from machine learning to meshless methods. *Acta Numerica*, 15:543–639, 2006.
- [8] B. Schölkopf and A.J. Smola. *Learning with Kernels*. MIT Press, Cambridge, 2002.
- [9] J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- [10] H. Wendland. *Scattered Data Approximation*. Cambridge University Press, 2005.

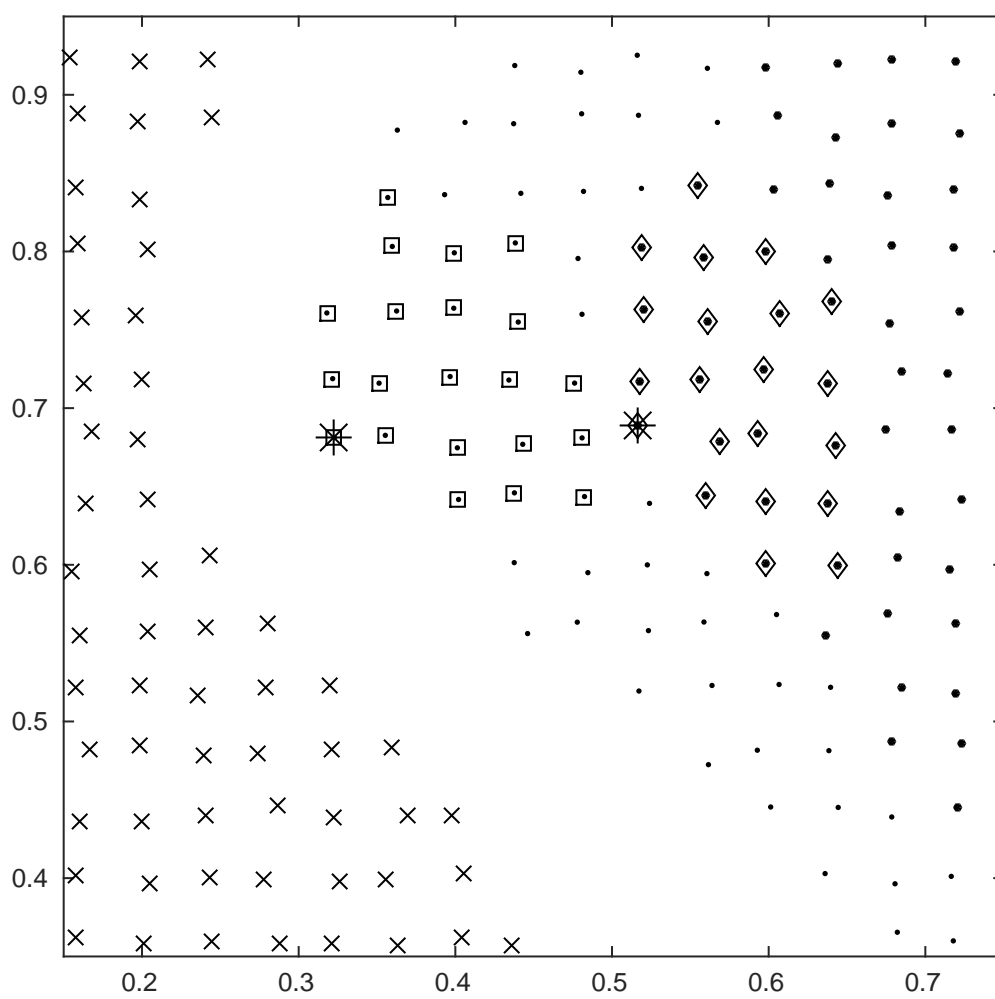


Figure 3: Localized blow-up, zoomed in

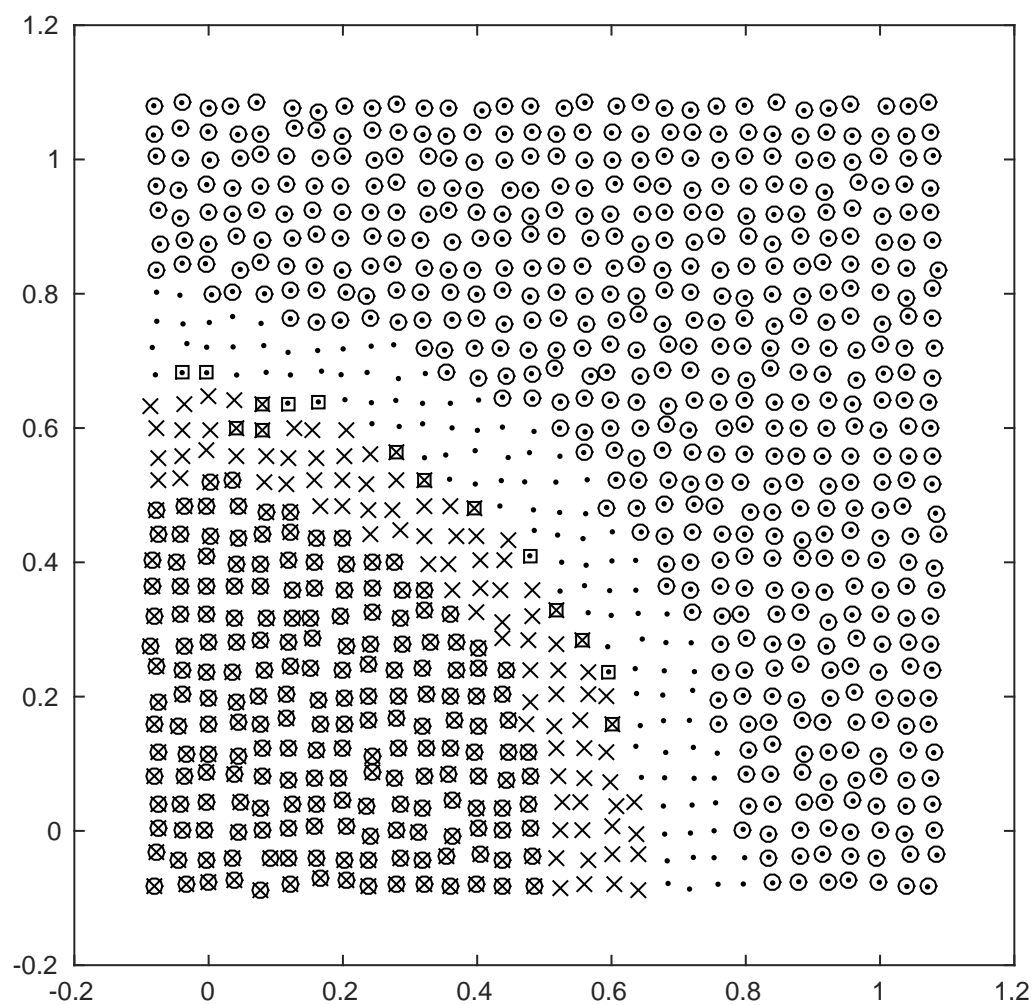


Figure 4: Example 3: class 1 as dots, class 2 as crosses

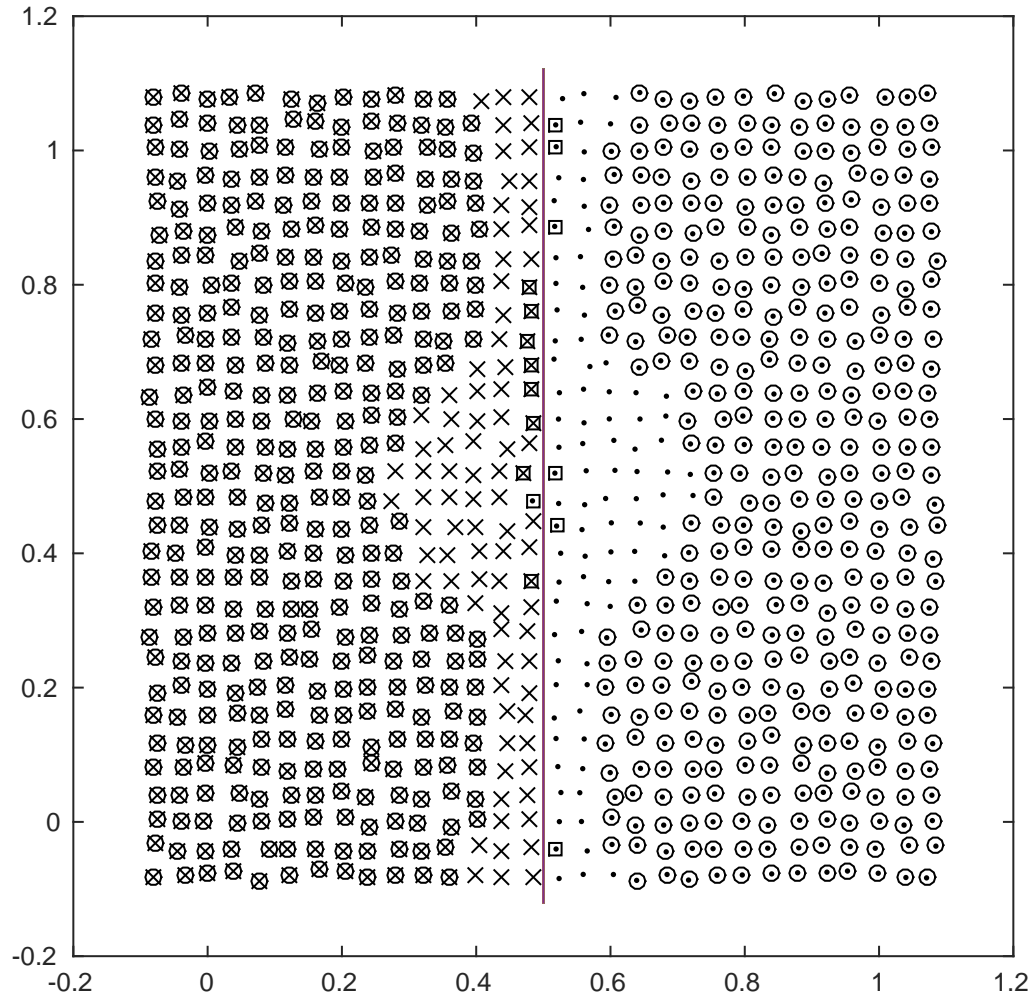


Figure 5: Example 4: class 1 as dots, class 2 as crosses